

CROSS-VALIDATION

We typically split data into D_{train} & D_{test} .

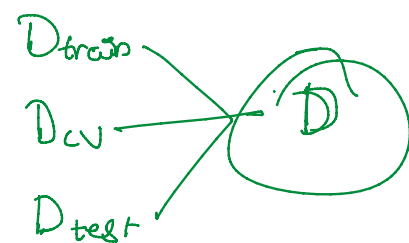
We train model on D_{train} & we get the accuracy on D_{test} .

But we can't generalize this result on future data. Because to evaluate / to do hyperparameter tuning we have used D_{test} here.

So we can't generalize the result on future unseen data.

So we need to have some other idea.

We can split our data into



D_{cv} → cross validation data.

D_{train} → Train the model.

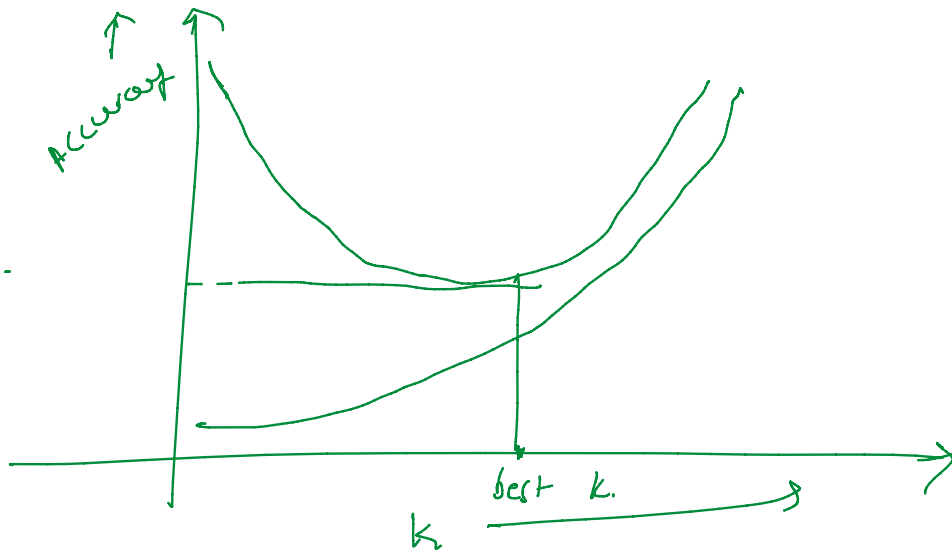
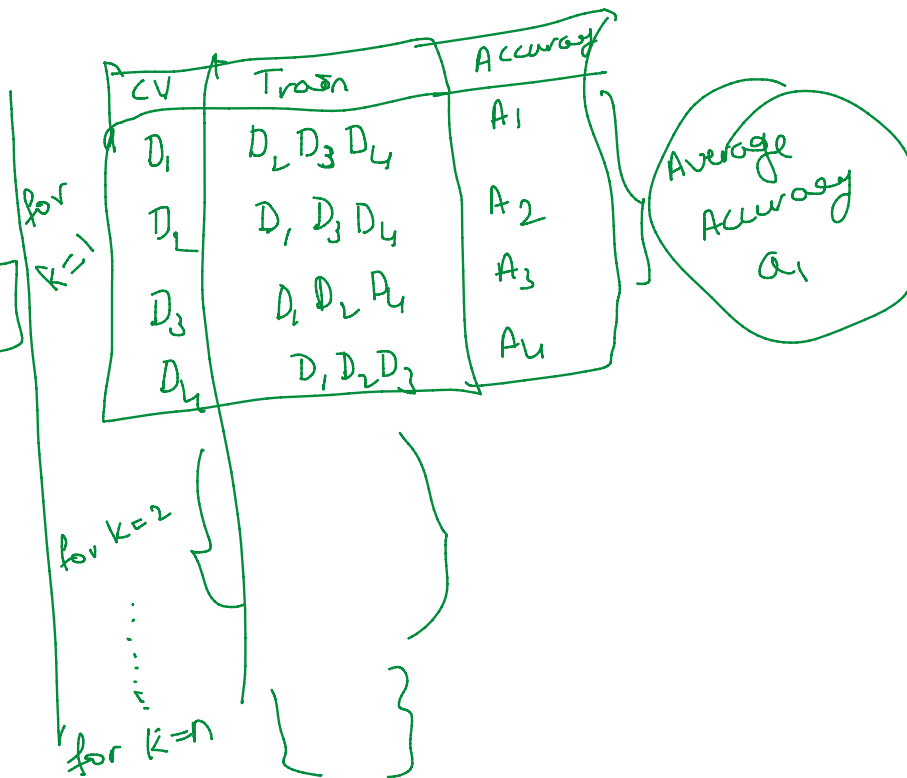
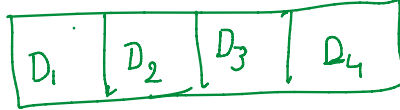
D_{cv} → hyperparameter tuning

D_{test} → Used as unseen data for generalization.

By splitting train Data into train & CV, we loose ... training. Hence for better performance

By splitting some data for training. Hence for better performance we use, k' -fold CV

typically $k'=10$



So, It is best practice to perform cross validation to find best hyperparameter.

or, — .
best hyperparameter.

